

Machine Learning

CSE 343/543

Lecture 3

Classification – Naïve Bayes, k-NN

Today's Lecture

- Classification
 - Probabilistic
 - Nonparametric
 - Nonmetric

Probabilistic Classifier: Bayesian

- For Binary Classification, we are interested in

$$p(y|\mathbf{x}) \quad y \in \{0, 1\}$$

- Applying Bayes Rule

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

- A Bayesian Classifier $\hat{y}_i = f(\mathbf{x}_i, \alpha)$ will try to match the training data distribution

$$p(f(\mathbf{x}_i, \alpha)|\mathbf{x}_i; \alpha) \sim p(y_i|\mathbf{x}_i; \alpha)$$

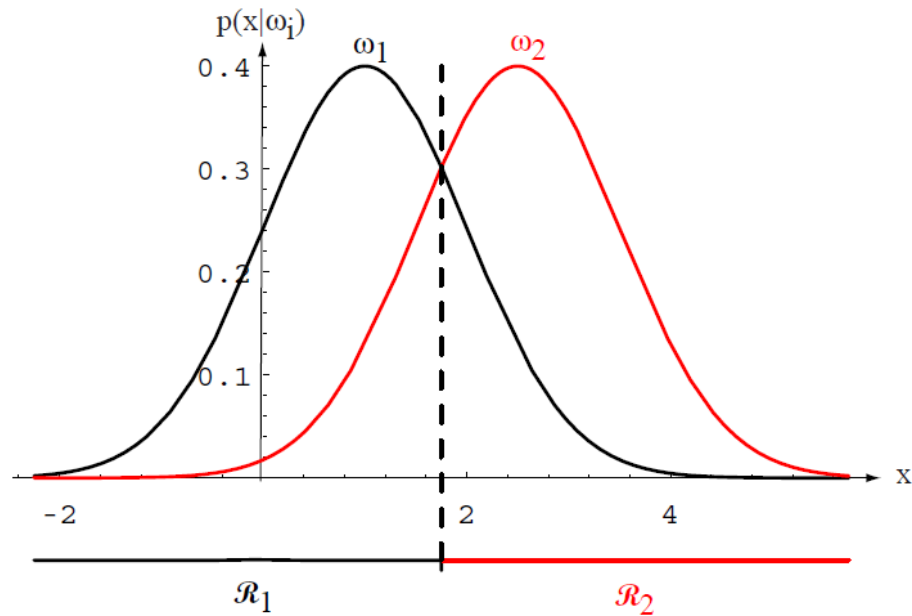
Probabilistic Classifier: Bayesian

- Classification is then performed as
 - If $p(f(\mathbf{x}_i, \alpha) = 1 | \mathbf{x}_i, \alpha) > p(f(\mathbf{x}_i, \alpha) = 0 | \mathbf{x}_i, \alpha)$
then the label
 $\hat{y}_i = 1$
else
 $\hat{y}_i = 0$

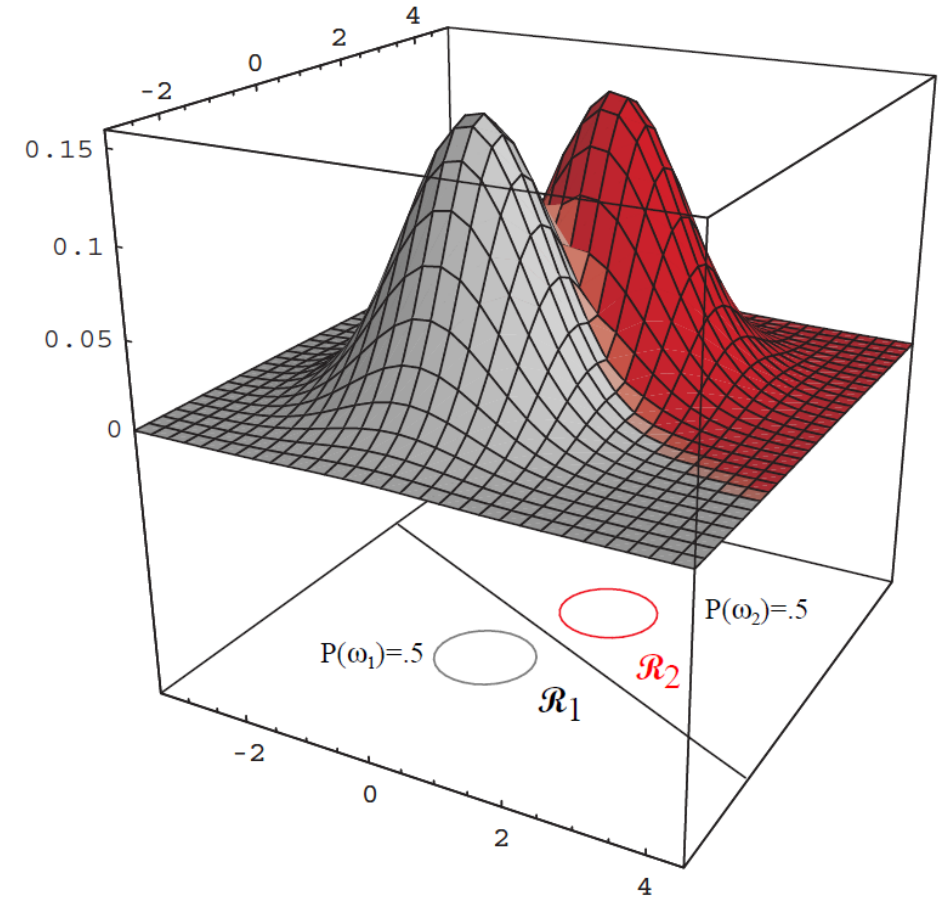
- For M-class classification

$$\hat{y}_i = \arg \max_{y=\{0,1,\dots,M\}} p(f(\mathbf{x}_i, \alpha) = y | \mathbf{x}_i)$$

Visualization: Bayes Classifier



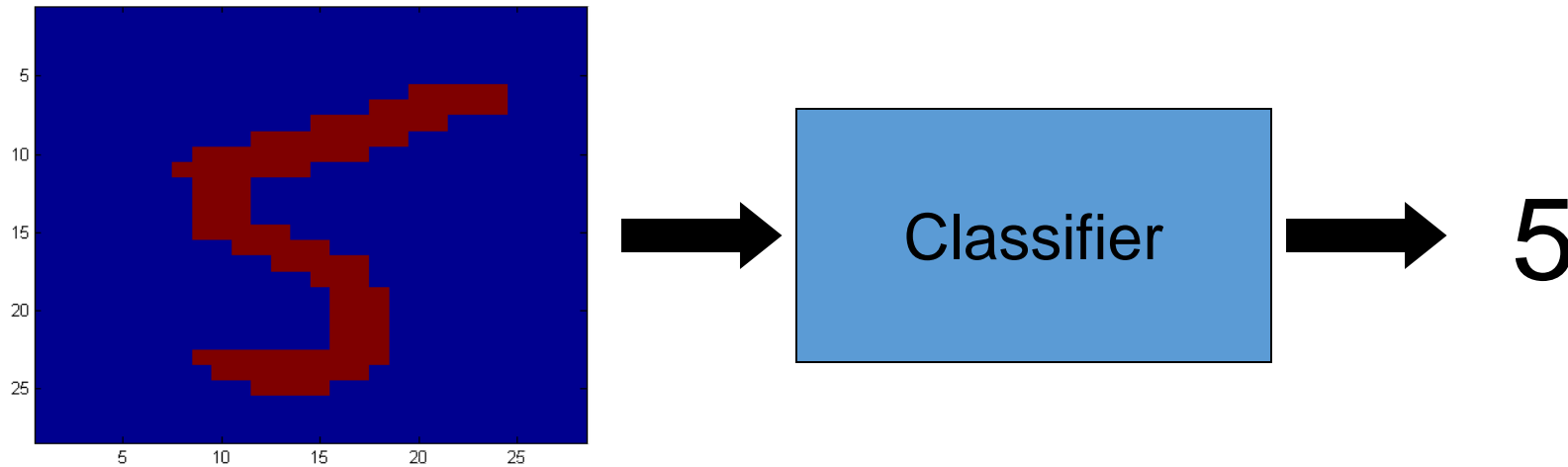
1-D



2-D

Concrete Example: Bayes Classifier

- **Digit Recognition**



- $X_1, \dots, X_n \in \{0,1\}$ (Black vs. White pixels)
- $Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

The Bayes Classifier

- We saw that a good strategy is to predict:

$$\arg \max_Y P(Y|X_1, \dots, X_n)$$

- (for example: what is the probability that the image represents a 5 given its pixels?)
- So ... How do we compute that?

The Bayes Classifier

- Uses Bayes Rule!

$$\begin{array}{c} \text{Posterior} \\ \downarrow \\ P(Y|X_1, \dots, X_n) \end{array} = \frac{\begin{array}{c} \text{Likelihood} \\ \downarrow \\ P(X_1, \dots, X_n|Y) \end{array} \begin{array}{c} \text{Prior} \\ \downarrow \\ P(Y) \end{array}}{\begin{array}{c} \uparrow \\ P(X_1, \dots, X_n) \\ \text{Normalization Constant} \end{array}}$$

- Why did this help?
 - Well... We think that we might be able to specify how features are “generated” by the class label

The Bayes Classifier

- Let's expand this for our digit recognition task:

$$P(Y = 5|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 5)P(Y = 5)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$
$$P(Y = 6|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 6)P(Y = 6)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$

- To classify, we'll simply compute these two probabilities and predict based on which one is greater

Model Parameters

- For the Bayes classifier, we need to “learn” two functions, the likelihood and the prior
- How many parameters are required to specify the prior for our digit recognition example?

Model Parameters

- How many parameters are required to specify the likelihood?
 - (Supposing that each image is 30x30 pixels)
- # of parameters for modeling $P(X_1, \dots, X_n | Y)$:
 - $2(2^n - 1)$
- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)

The Naïve Bayes Model

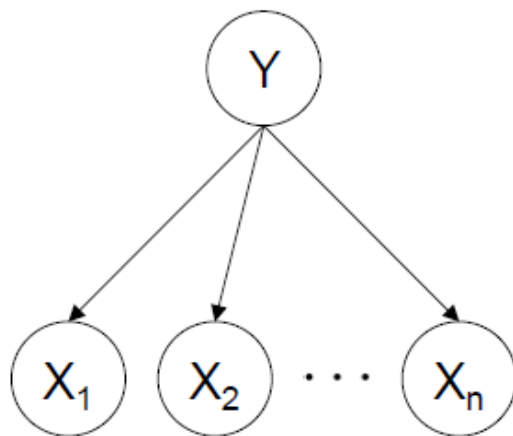
- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

- # of parameters for modeling $P(X_1 | Y), \dots, P(X_n | Y)$
 - $2n$

The Naïve Bayes Classifier

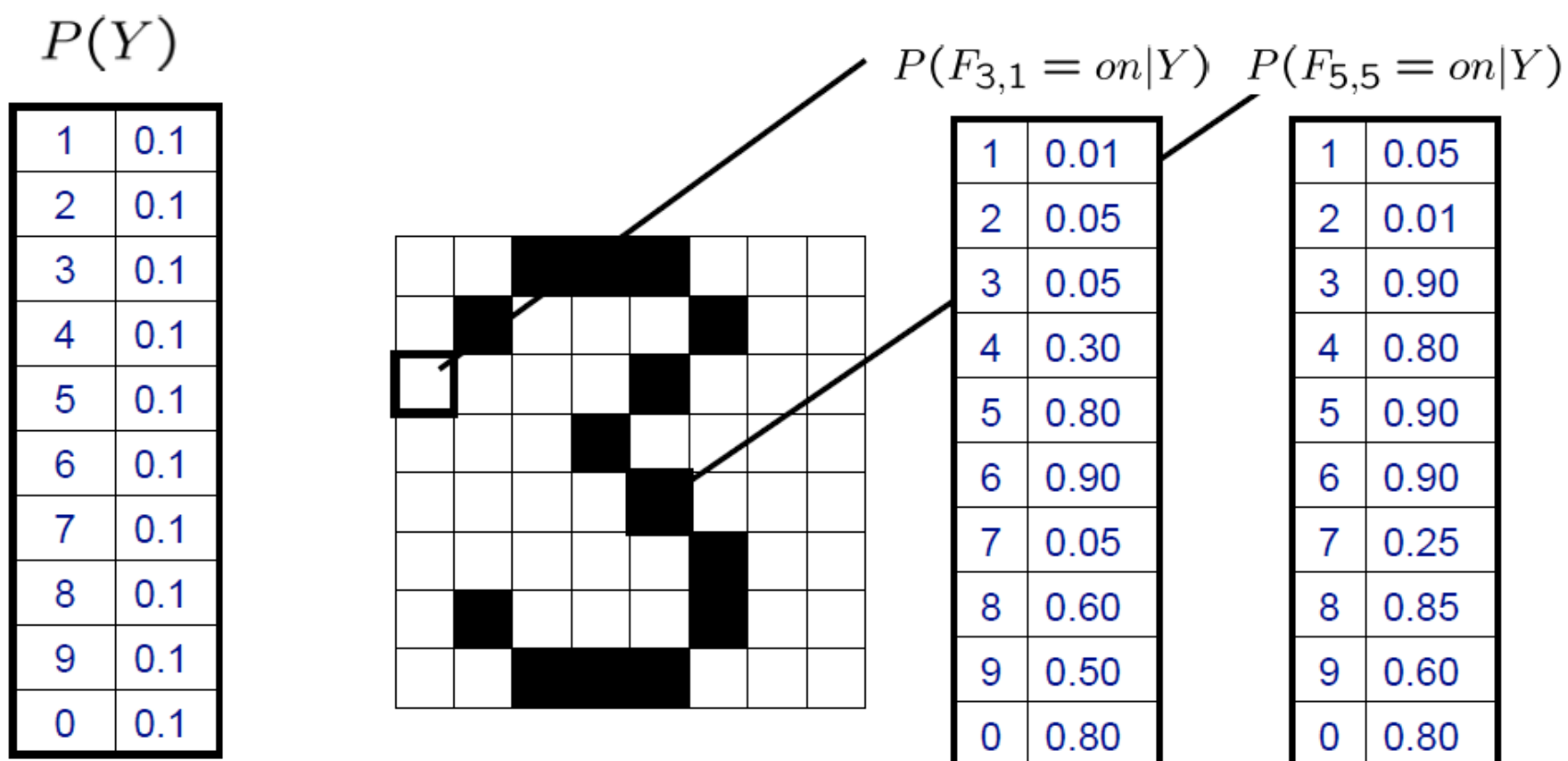
- Given:
 - Prior $P(Y)$
 - n conditionally independent features \mathbf{X} given the class Y
 - For each X_i , we have likelihood $P(X_i | Y)$



- Decision rule:

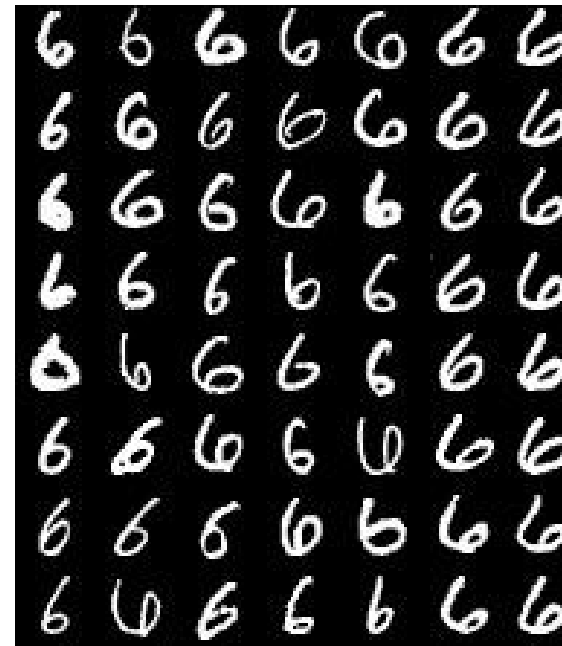
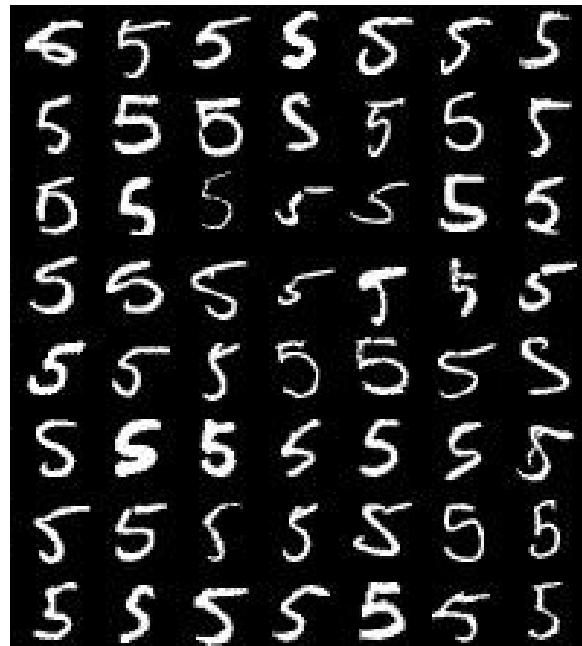
$$\begin{aligned}\hat{y} &= \arg \max_y P(y) P(x_1, \dots, x_n | y) \\ &= \arg \max_y P(y) \prod_i P(x_i | y)\end{aligned}$$

What has to be learned?



Naïve Bayes Training

- Now that we've decided to use a Naïve Bayes classifier, we need to train it with some data:



MNIST Training Data

Naïve Bayes Training

- Training in Naïve Bayes is **easy**:
 - Estimate $P(Y=v)$ as the fraction of records with $Y=v$

$$P(Y = v) = \frac{\text{Count}(Y = v)}{\# \text{ records}}$$

- Estimate $P(X_i=u | Y=v)$ as the fraction of records with $Y=v$ for which $X_i=u$

$$P(X_i = u | Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v)}{\text{Count}(Y = v)}$$

- (This corresponds to Maximum Likelihood estimation of model parameters)

Naïve Bayes Training

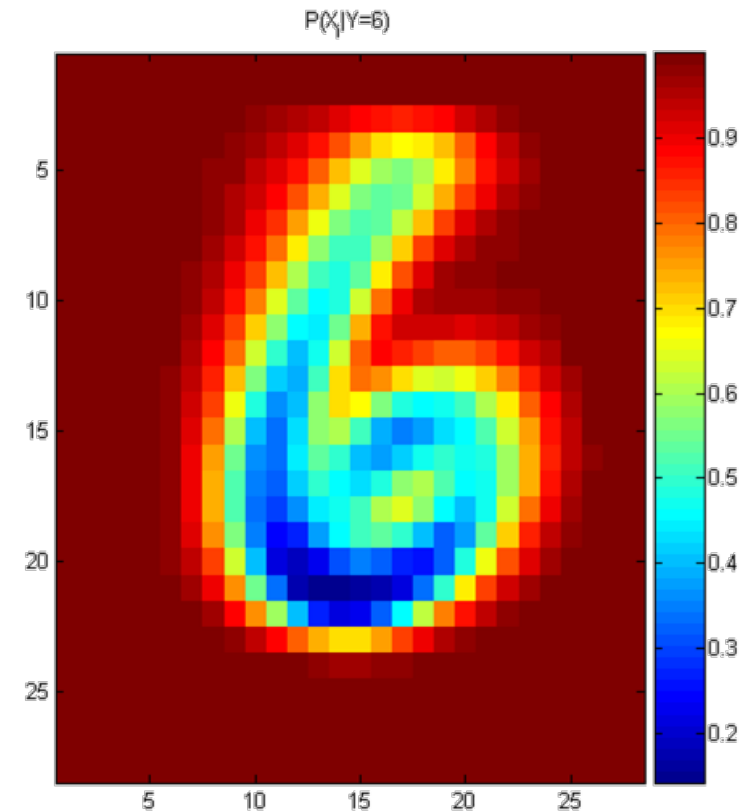
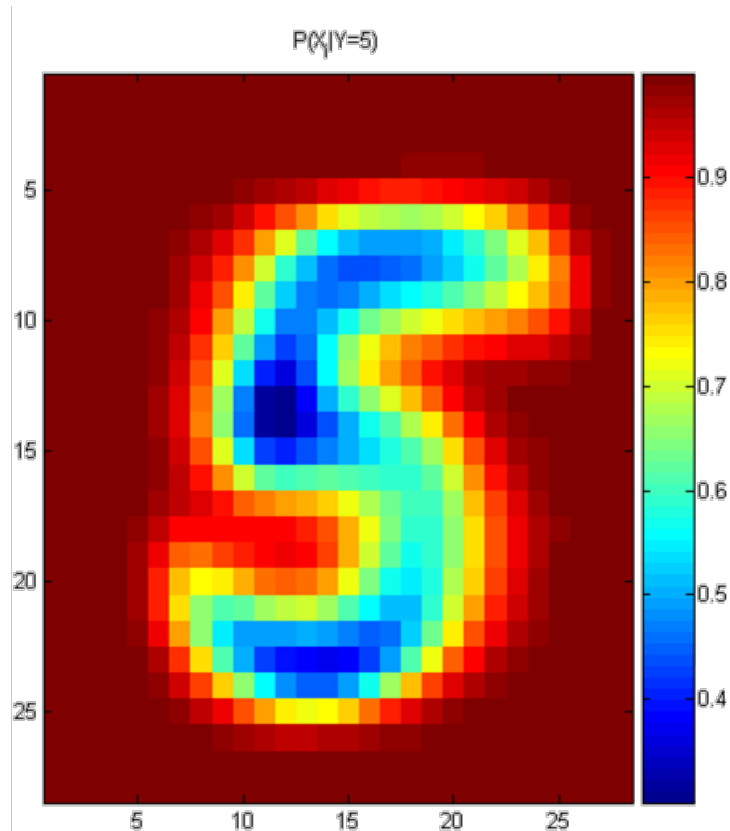
- In practice, some of these counts can be zero
- Fix this by adding “virtual” counts:

$$P(X_i = u|Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v) + 1}{\text{Count}(Y = v) + 2}$$

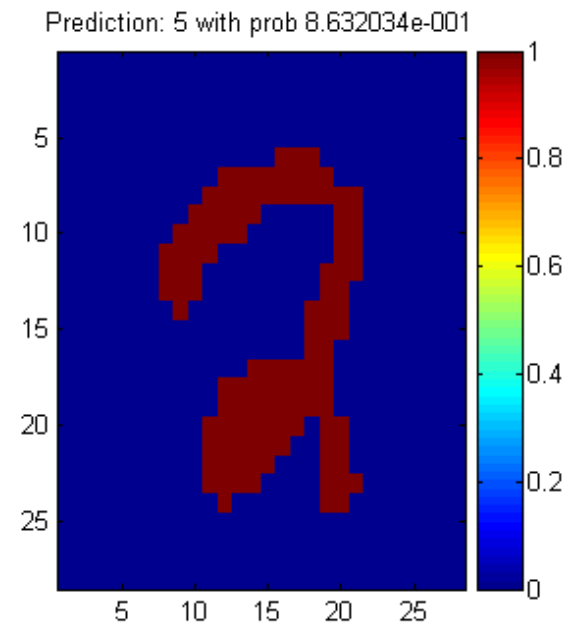
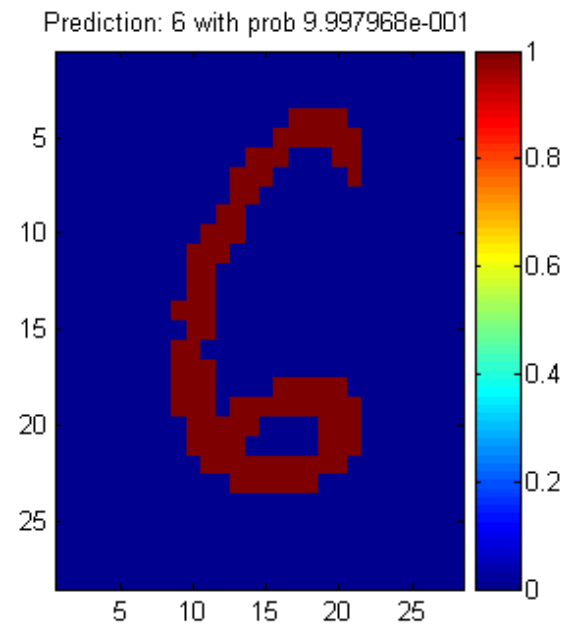
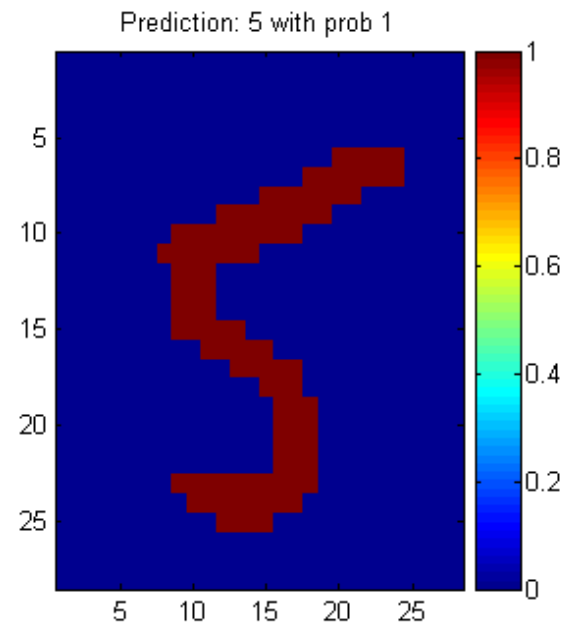
- (This is like putting a prior on parameters and doing MAP estimation instead of MLE)
- This is called *Smoothing*

Naïve Bayes Training

- For binary digits, training amounts to averaging all of the training fives together and all of the training sixes together.

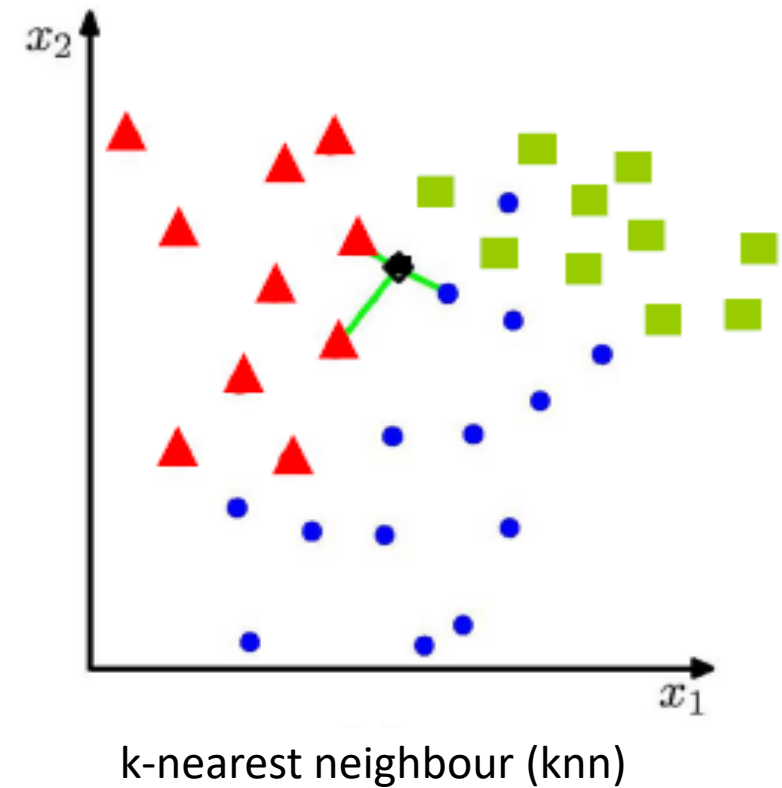


Naïve Bayes Classification



Non-parametric Classification: K-Nearest Neighbors

- Collect training data (x,y)
- For test sample
 - Find nearest neighbors in training set
 - Assign class label based on some consensus
- Parameters
 - Distance Metric for 'nearest' neighbors
 - Voting strategy



Non-metric Classification: Decision Trees

